

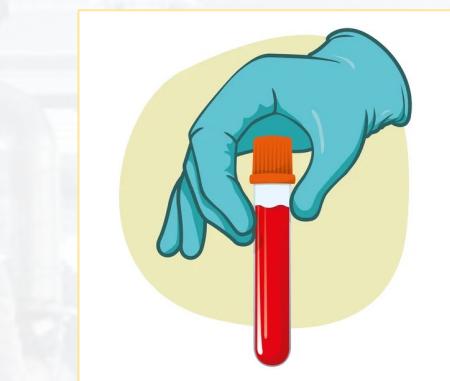
# Flexible pipeline for identification reliable genetic markers in large genomic datasets



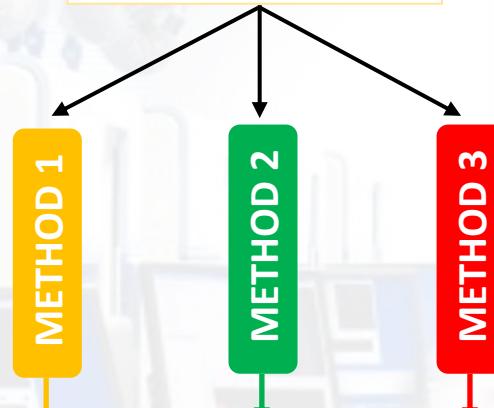
Gap.Jumper

Pawel Rosikiewicz  
Data Scientist, Pipeline Developer

## False Positive Leakage in BIG DATA



EVEN IF WE USE  
ONE BLOOD  
SAMPLE



WE WILL STILL GET  
DIFFERENT RESULTS  
FROM TIME TO TIME  
(..... ATGC....)



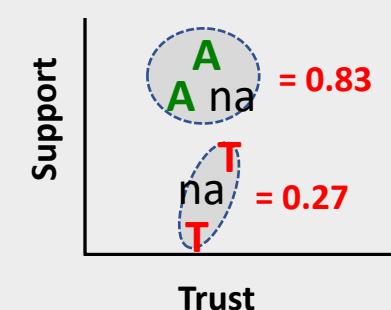
It becomes a major  
problem when:

- 👉 DETECTING RARE MUTATIONS
- 👉 USE LOWER QUALITY DATA
- 👉 CAN NOT APPLY IMPUTATION

## My Solution

I IMPLEMENTED APPROACH USED  
FOR NAVIGATION IN DRONES & AIRPLANES

GapJumper  
calculates  
Evidence Scores  
based only on  
empirical data



- New algorithm
- ✓ Kalman Filter
  - ✓ Semi-Supervised
  - ✓ DST-Based
  - ✓ No Imputation

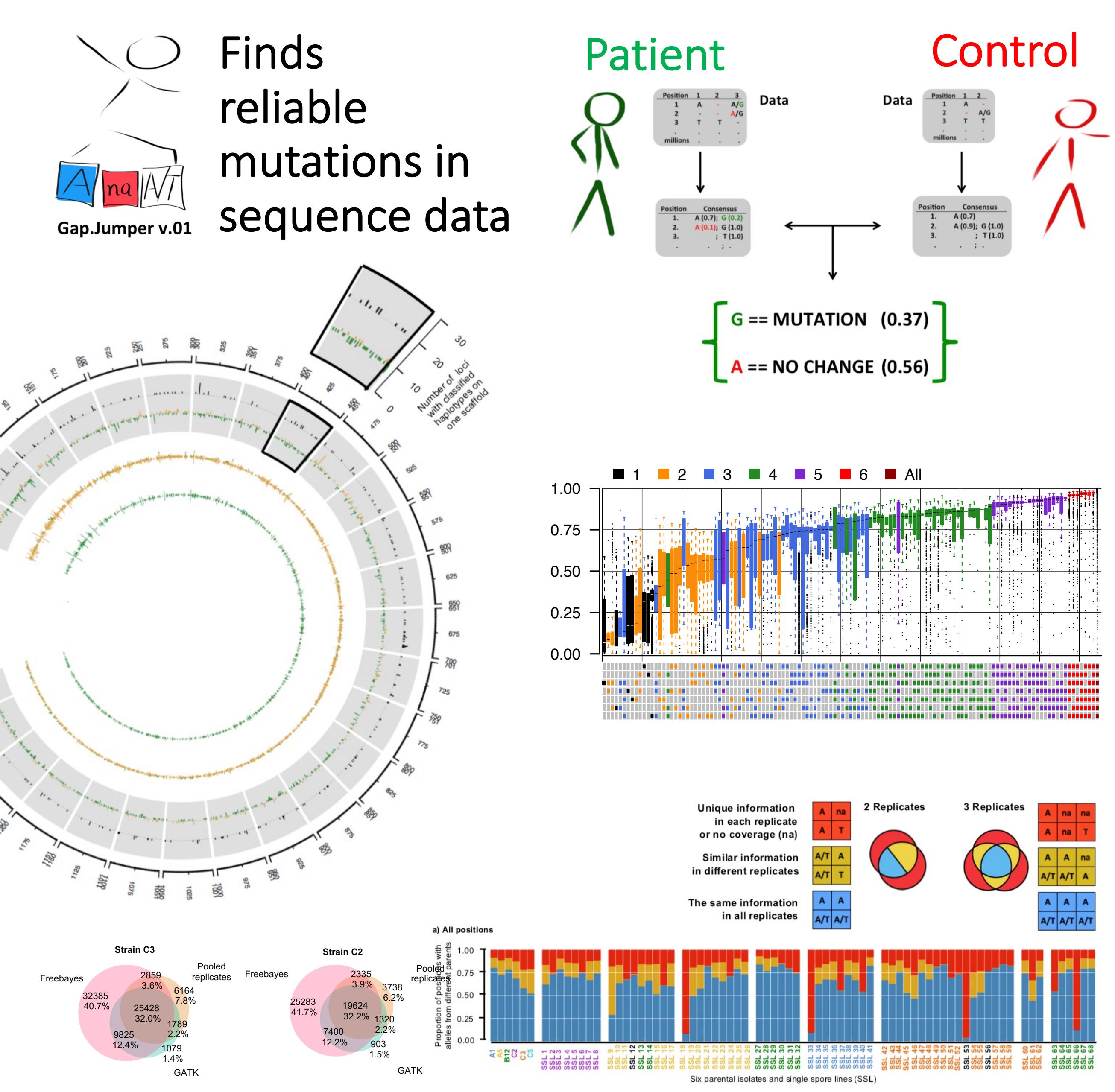
## NEW POSSIBILITIES FOR

- 👉 SCREENING STUDIES WITH LOW-COST SEQUENCING METHODS
- 👉 QC & SEQUENCE DATA INTEGRATION

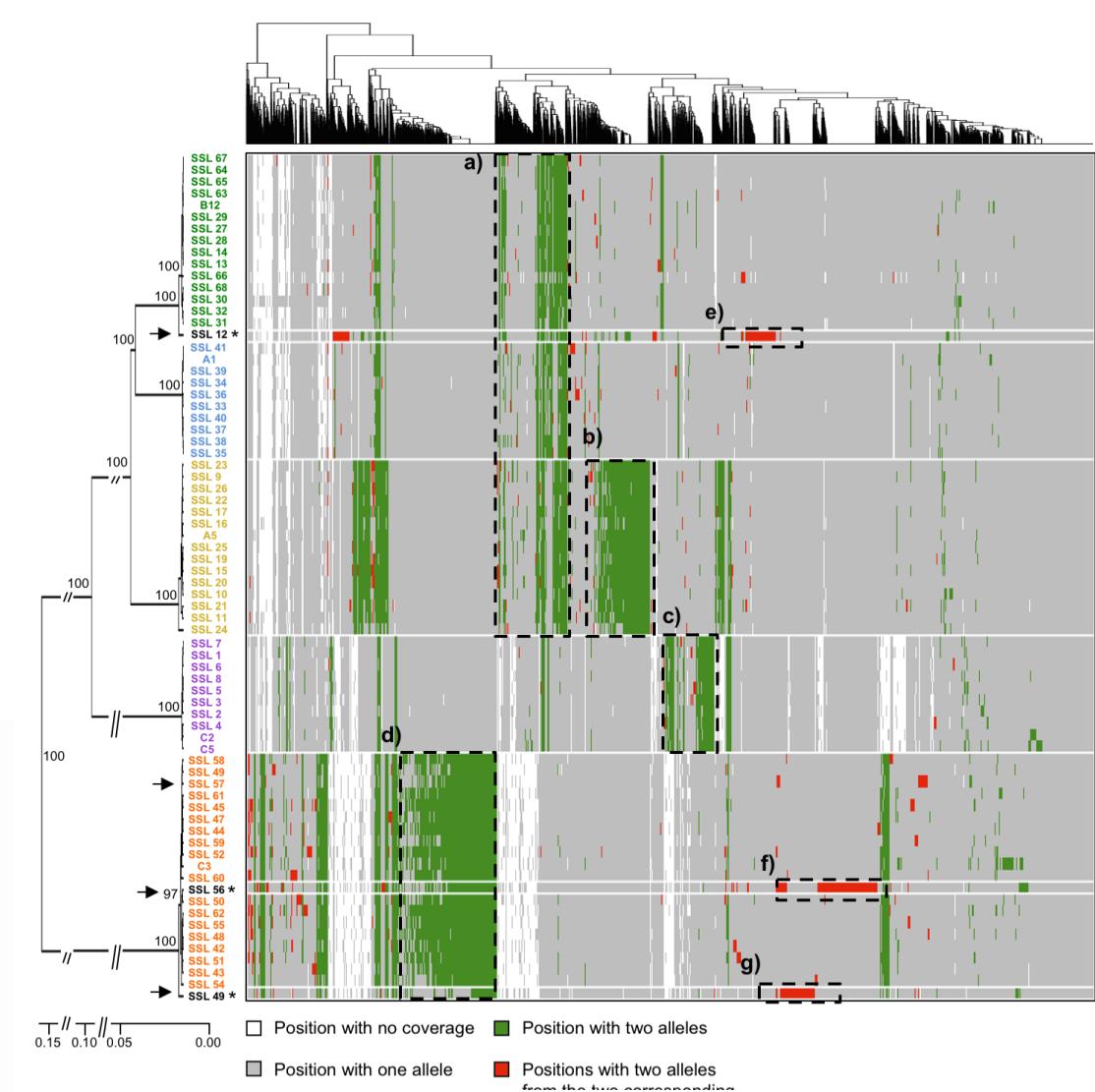
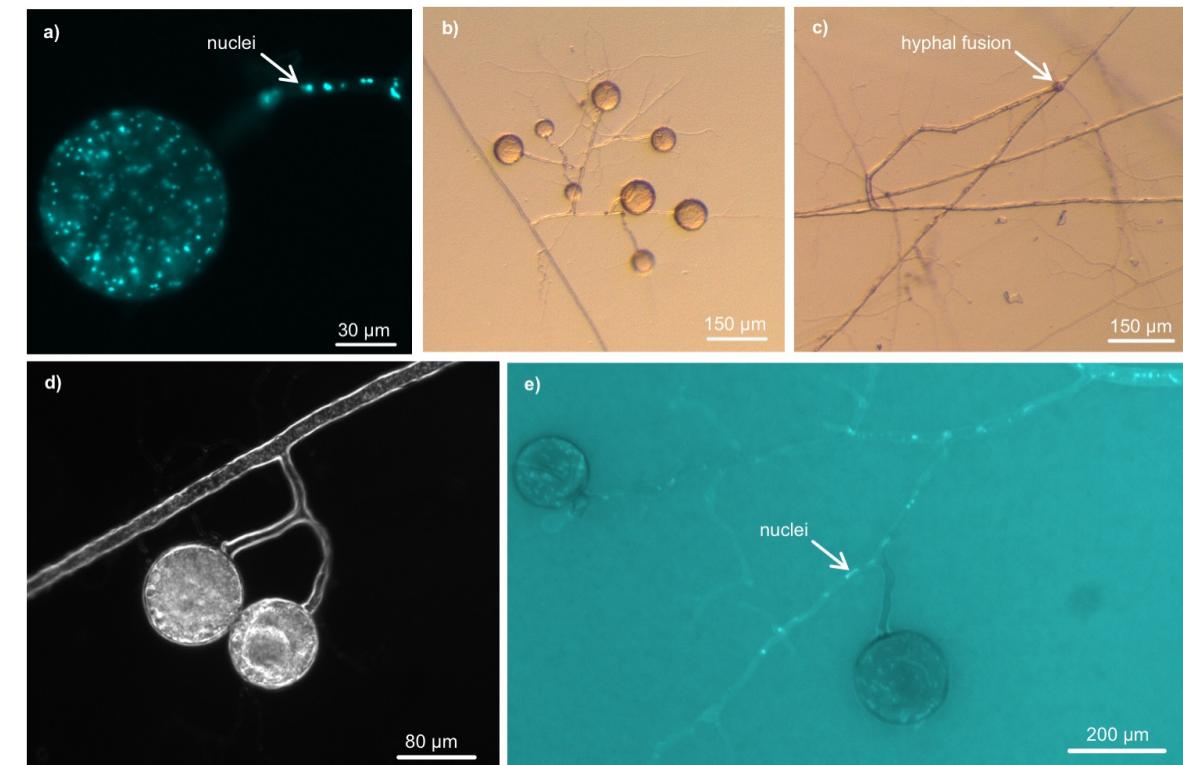
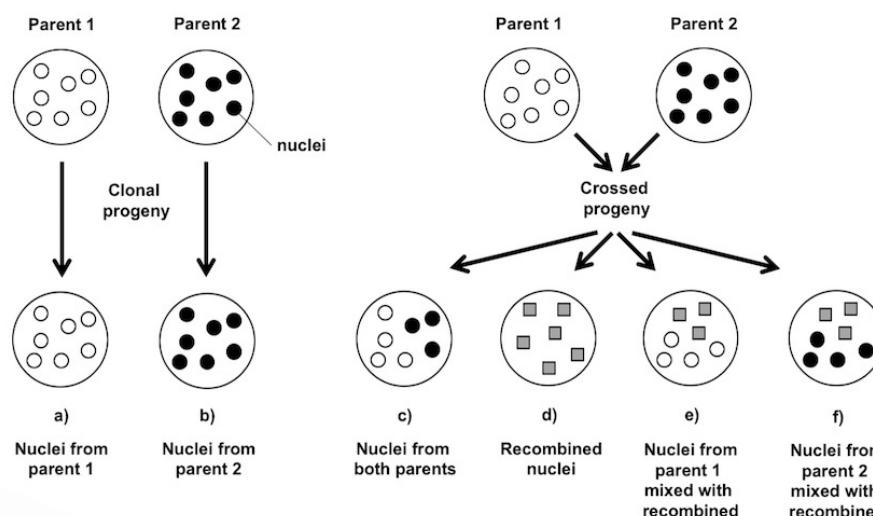
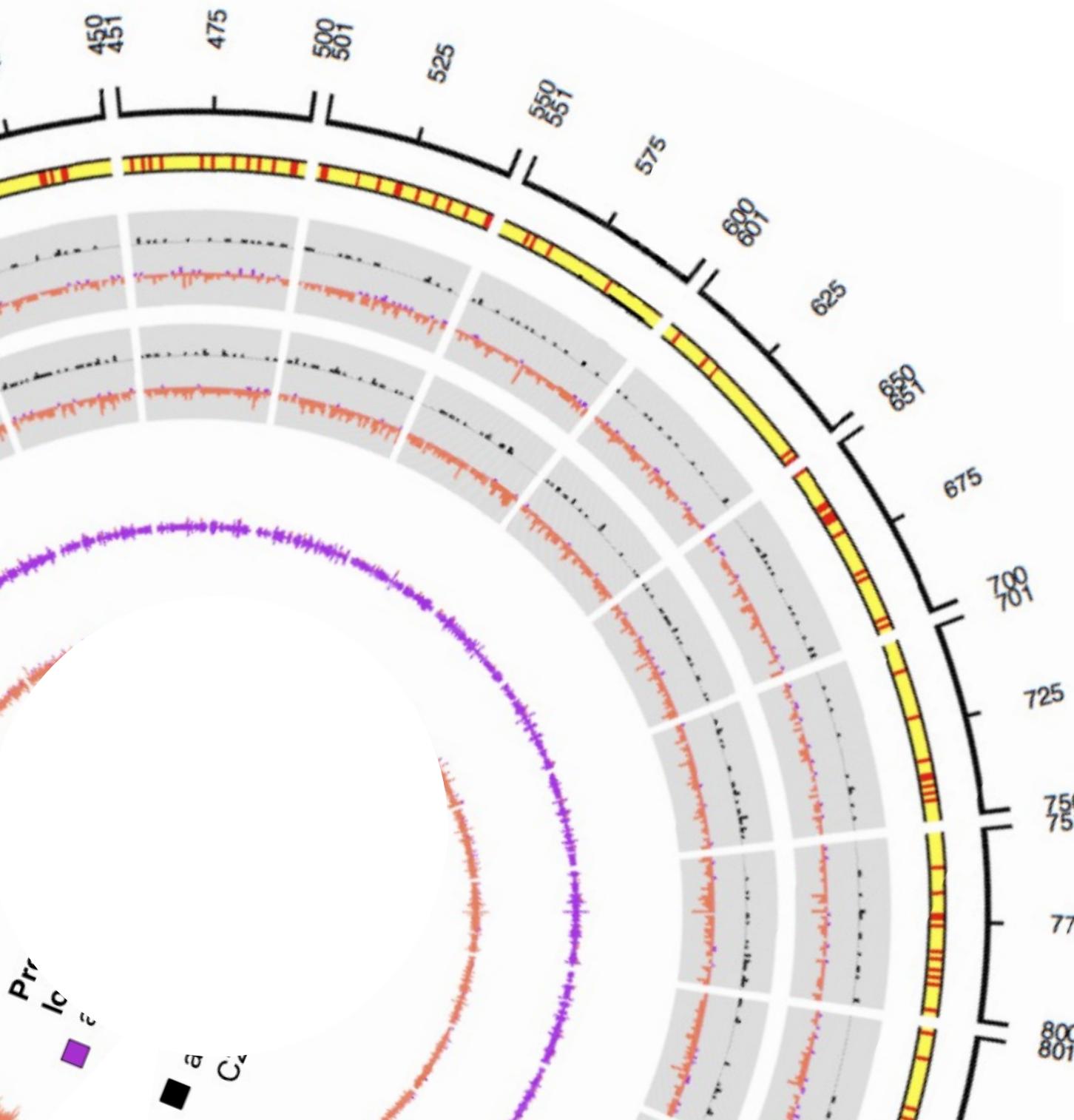
open source

## GapJumper on GitHub

<https://github.com/PawelRosikiewicz/DeepGenome>

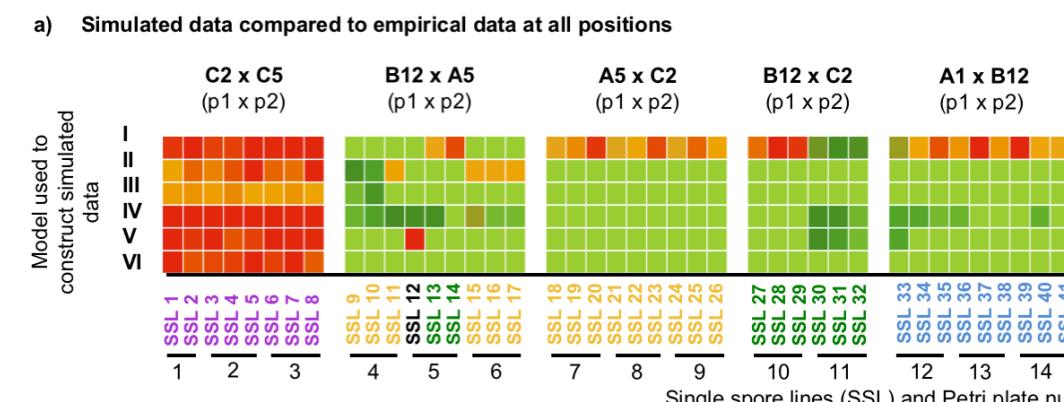


# Production of new Strains Biofertilizer Fungi



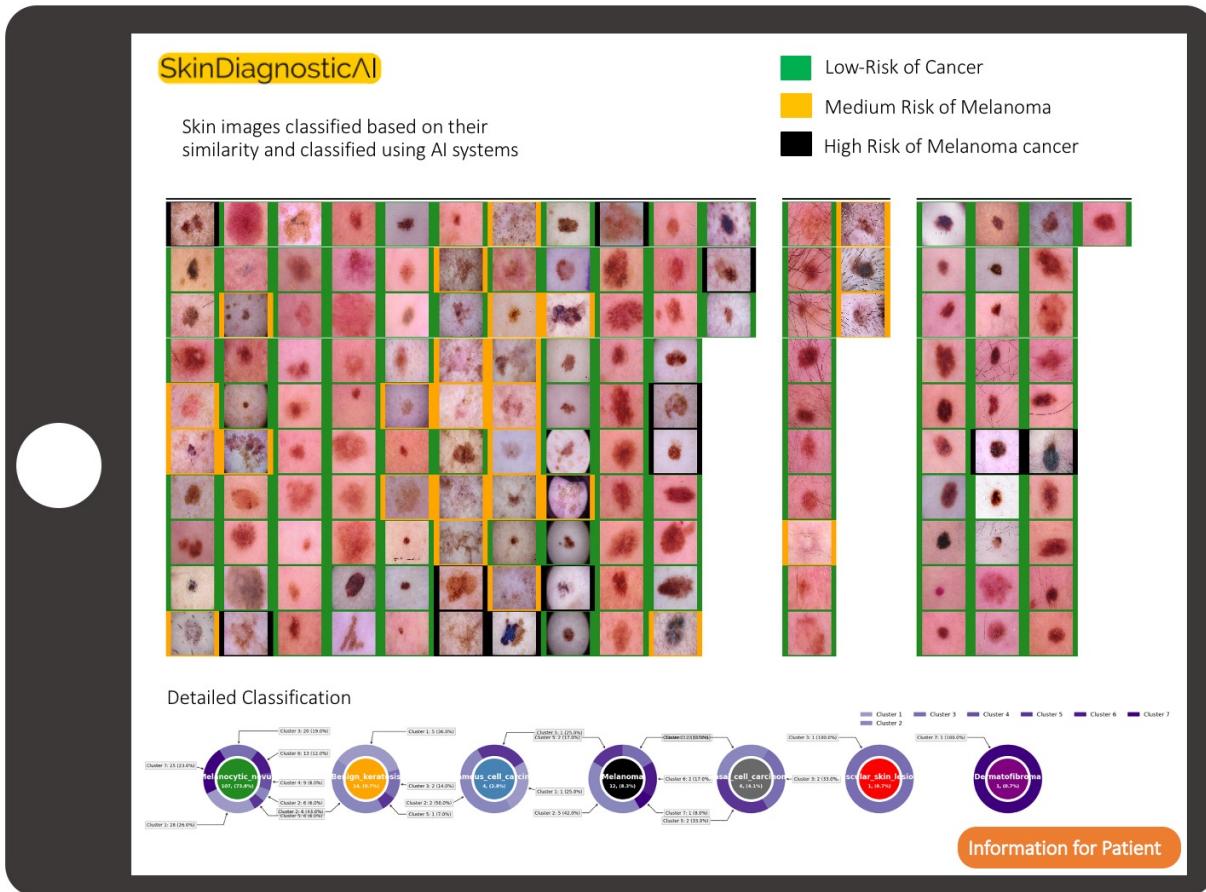
**Six different models used to construct simulated data (I-VI)**

Clonal lines	Crossed lines
I Clone of parent 1 (p1)	III Admixture of nuclei from both parents
II Clone of parent 2 (p2)	IV Recombined nuclei
	V Recombined nuclei mixed with nuclei from parent 1
	VI Recombined nuclei mixed with nuclei from parent 2



Single spore lines (SSL) and Petri plate numbers

## MY ROLE: Product Owner, Data Scientist



### Technical goals

#### Define MVP requirements

What product to build?

#### Deploy it on the cloud

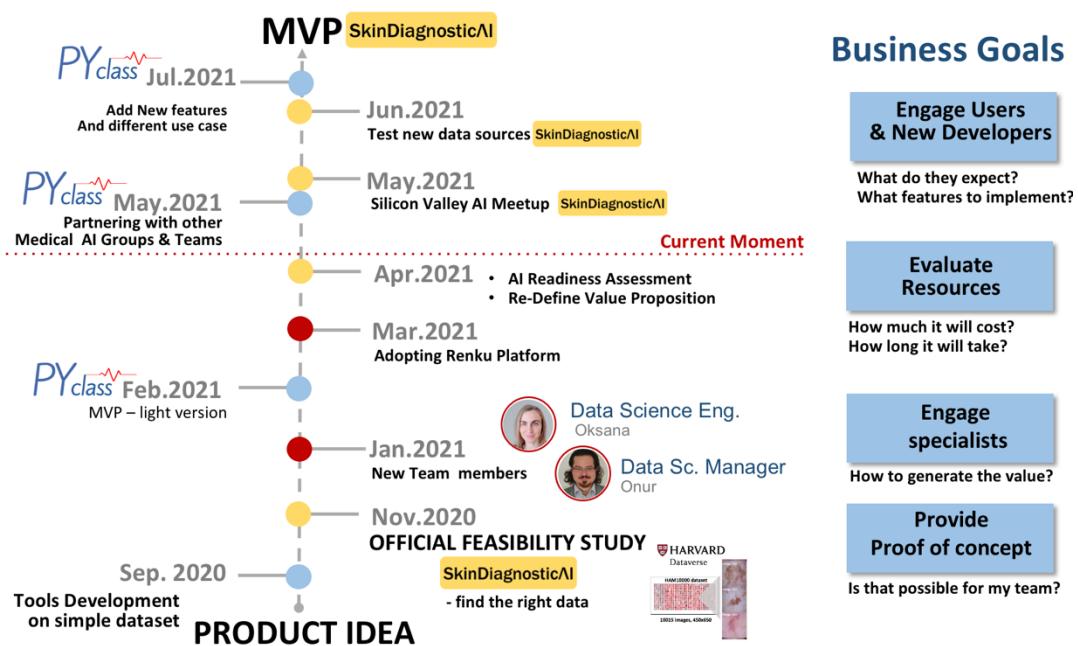
What model & data to use?

#### Improve the design

What features are important?

#### Build The Product

How to do it fast?



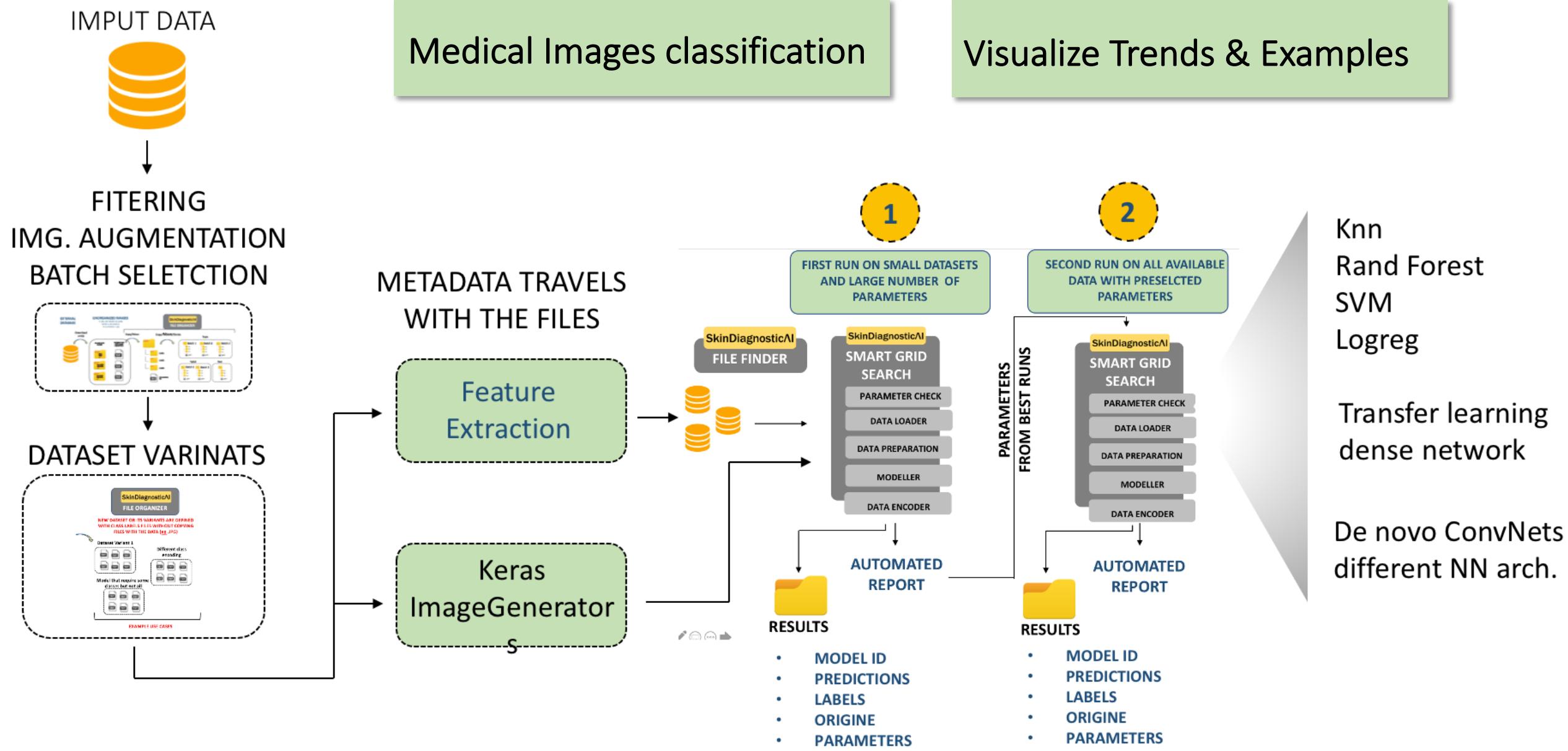
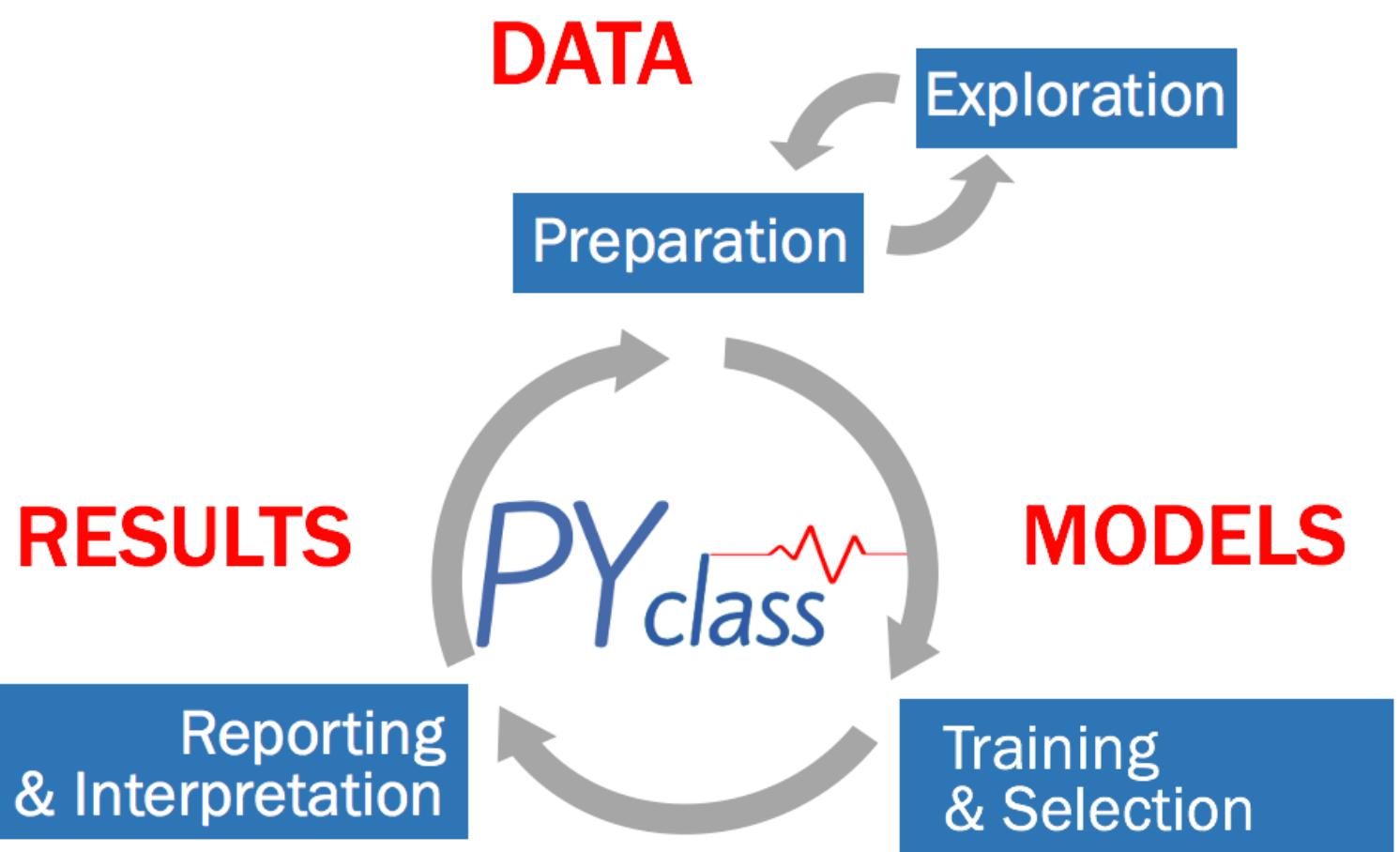
### LINKS TO MY WORK

- ✓ Project Website: <https://simpleai.ch/skin-diagnostic-ai/>
- ✓ GitHub: <https://github.com/PawelRosikiewicz/SkinDiagnosticAI>
- ✓ My Talk on Meetup in Silicon Valley: <https://youtu.be/W624gdkDqRQ?t=1259>



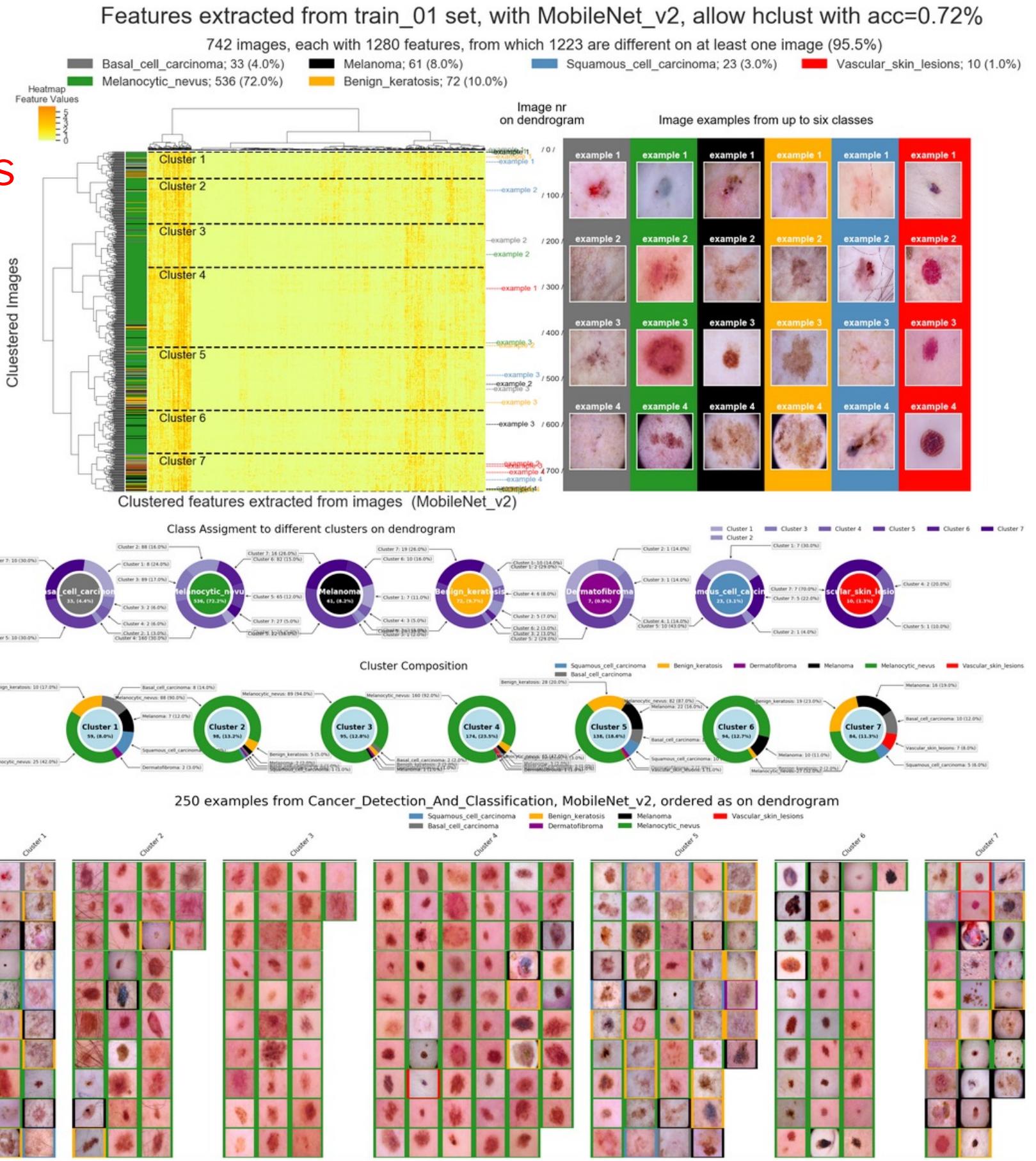
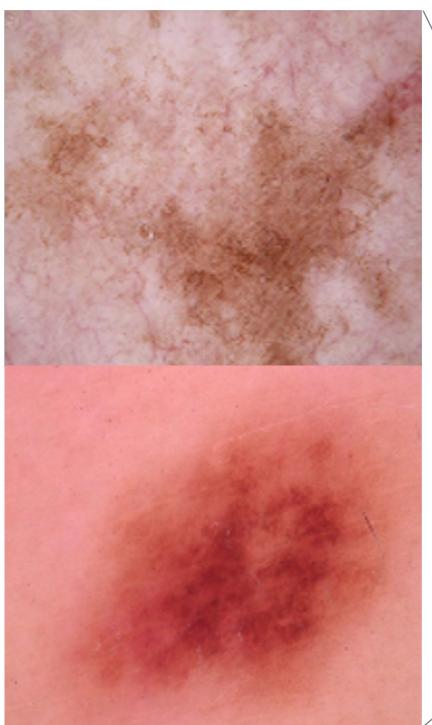
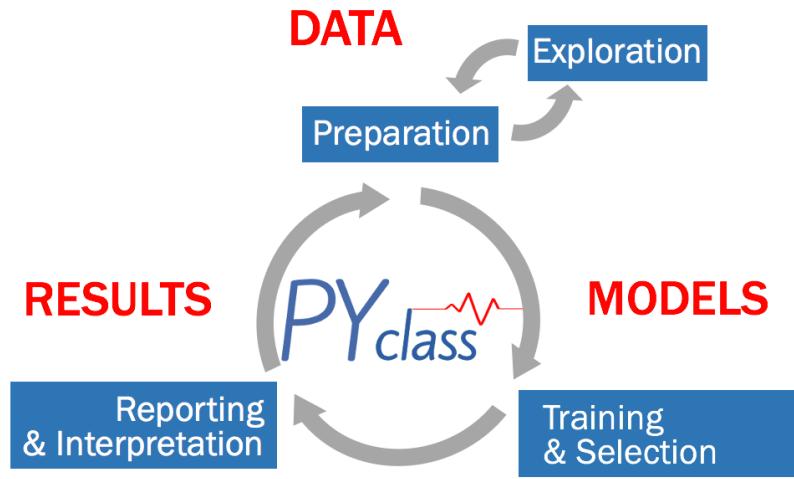
AI workbench for Analysis  
of Medical Images

Easy interpretation





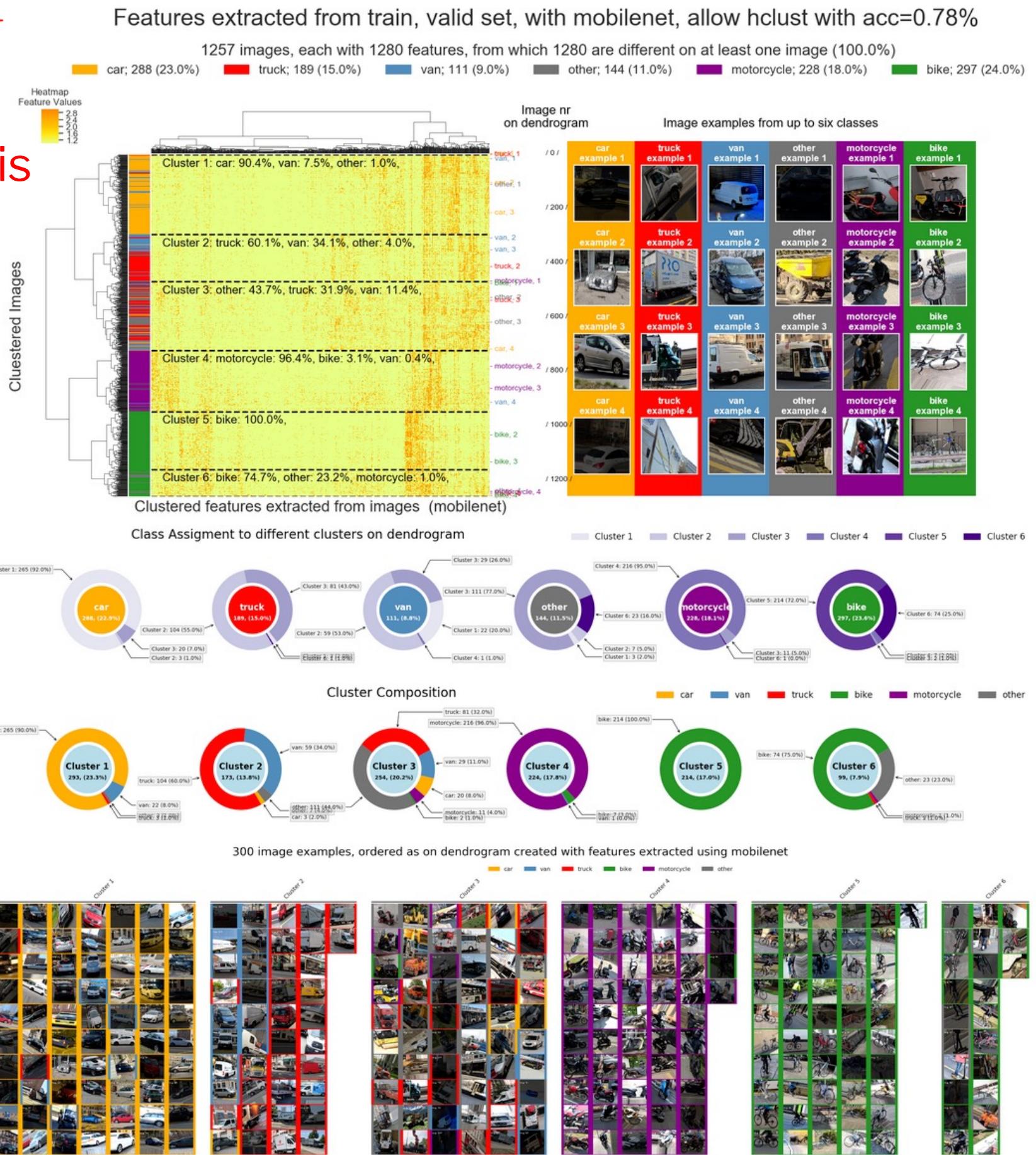
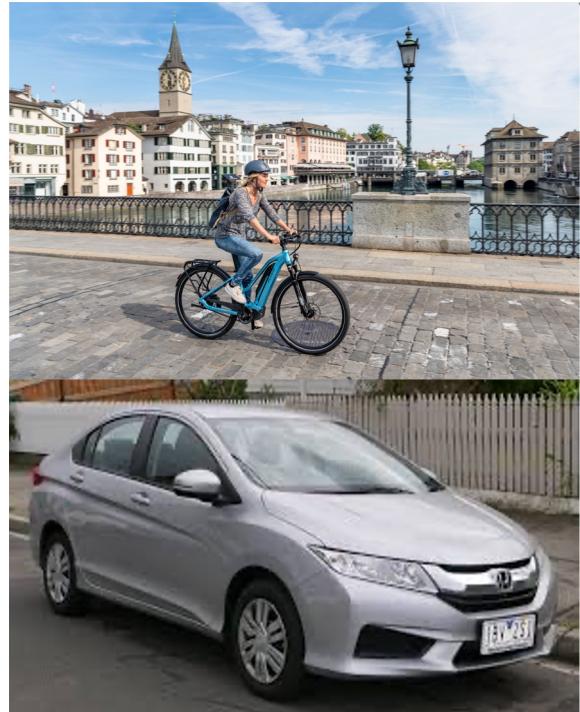
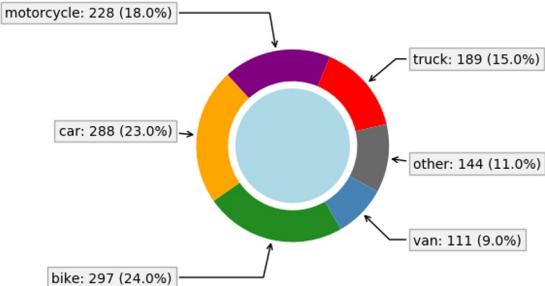
# AI workbench for Analysis of Medical Images



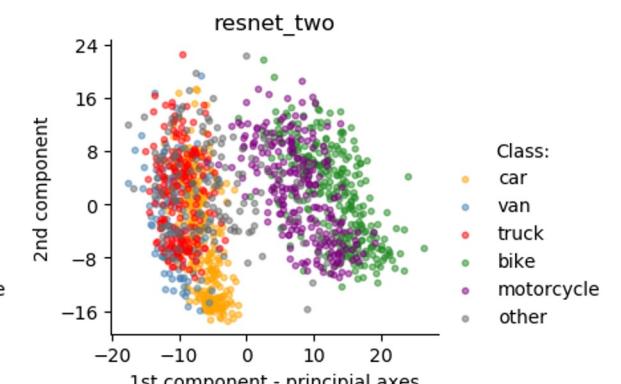
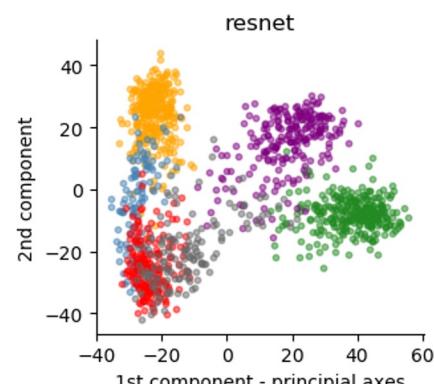
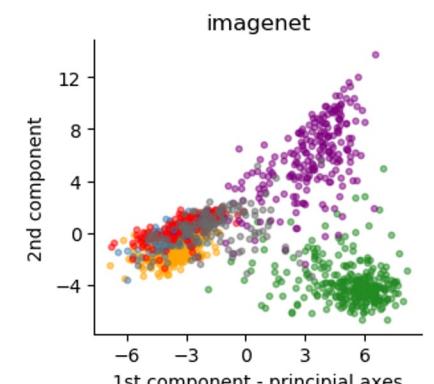
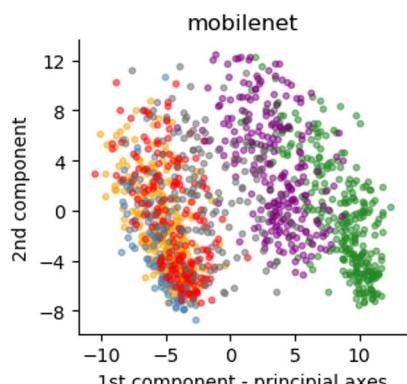


# AI workbench for Analysis of Medical Images

Swissroads Dataset Composition  
- including augmented images -



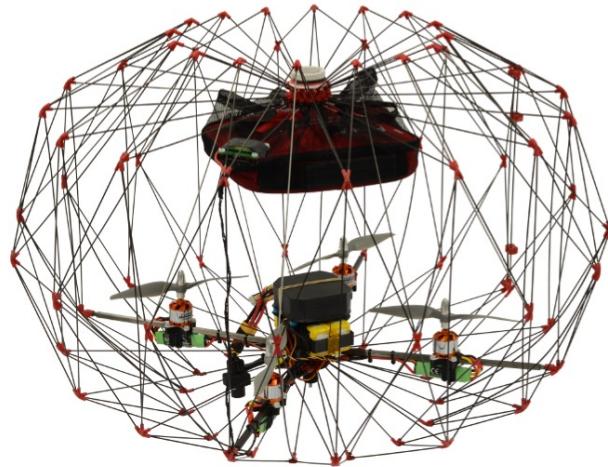
PCA: visualization of extracted features using the first two principal components



# Project Manager @ DRONISTICS



## Medical Delivery Drones



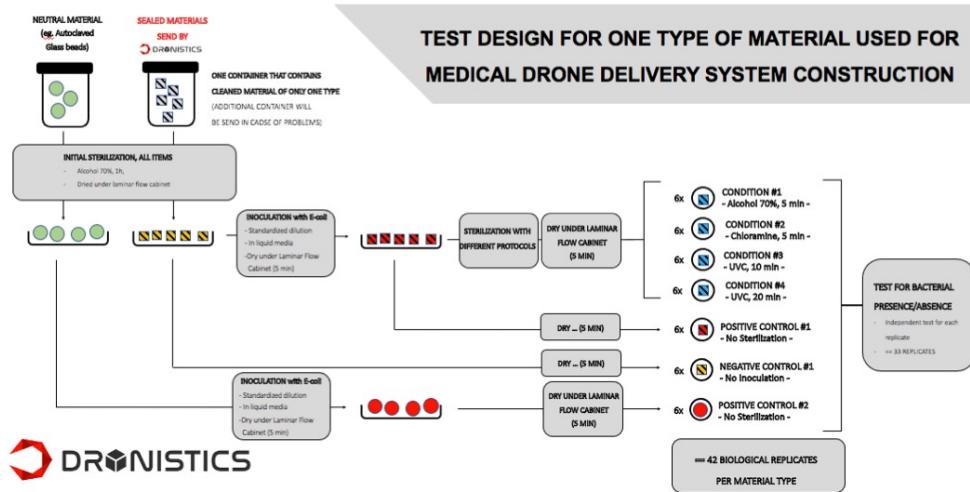
### Experiment Design

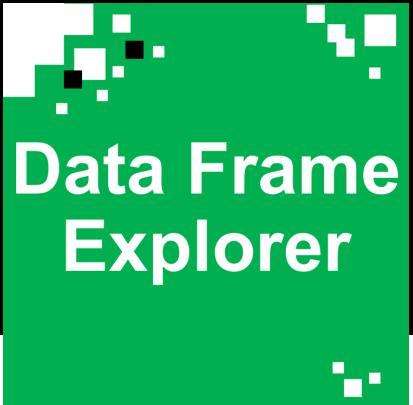
### Project & Team Management

### Data Analysis

## MY WORK

I led R&D team @ King's College London on adapting drone maintenance procedures to requirement from hospitals, & rescue services.

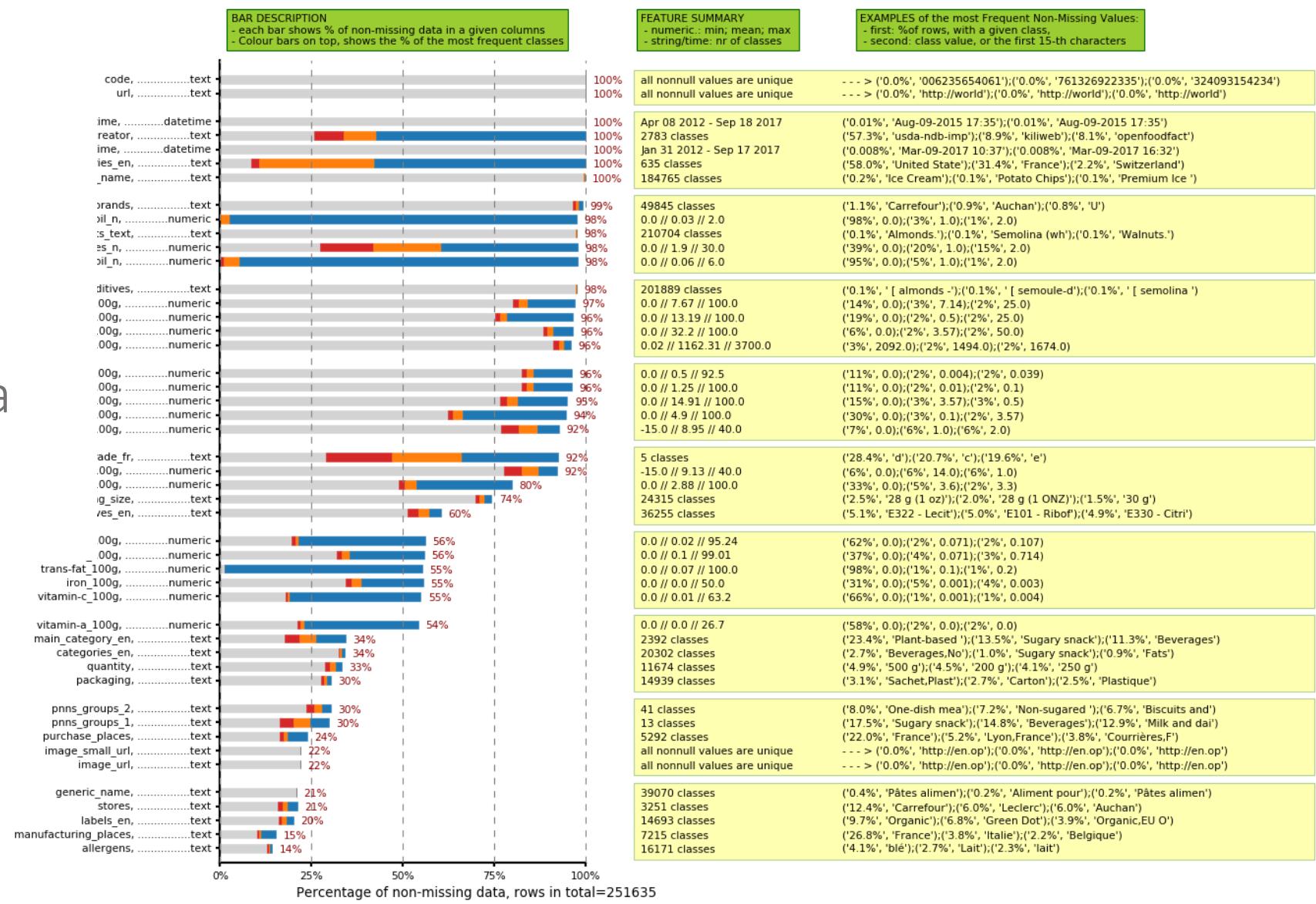




# My open-source library for fast exploration, analysis, and cleaning of Python Data Frames

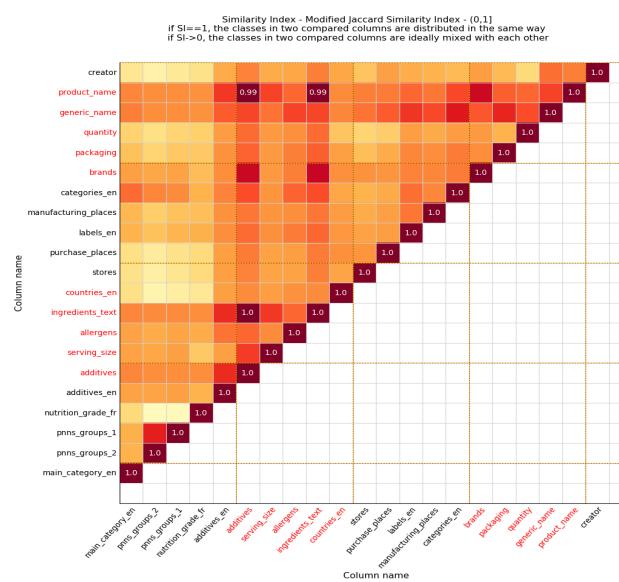
## Fast summary for hundreds of features

- ✓ Data type
- ✓ Value examples
- ✓ Amount of missing data
- ✓ Value ranges
- ✓ Class number

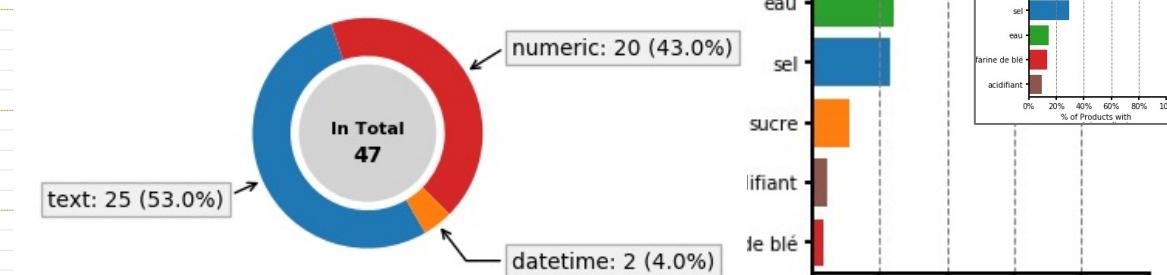


## Helps finding

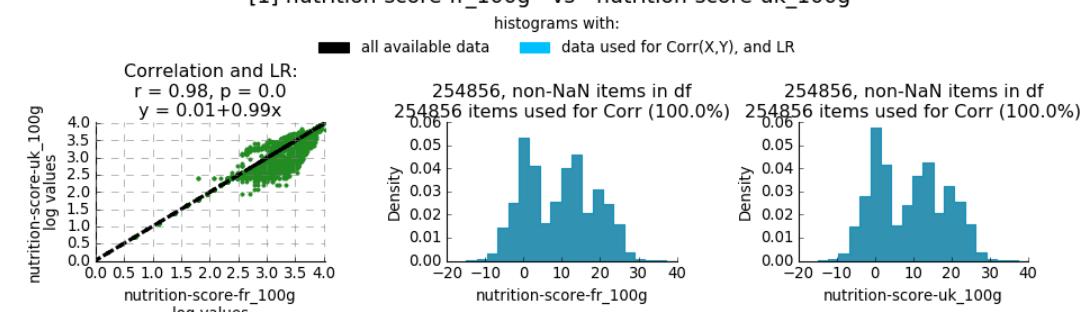
- ✓ Duplicates
- ✓ Correlations
- ✓ Cleaning and organizing the data
- ✓ Explore each feature



Column number with each data type



[1] nutrition-score-fr\_100g - vs - nutrition-score-uk\_100g



## USE CASE EXAMPLE:

<https://github.com/PawelRosikiewicz/HealthyFoodLabels>





# Paweł Rosikiewicz

## Lead Researcher on end-to-end multi-stage project

*Created Entire System for production & identification of polyploid microorganisms*

LABOLATORY

MOLECULAR BIOLOGY  
& BIOINFORMATICS

NGS  
DATA ANALYSIS

MY ROLE



DESIGN & DEV.



NGS SEQUENCING



GENETICS



MANAGED LAB TEAM



VARIANT CALLING

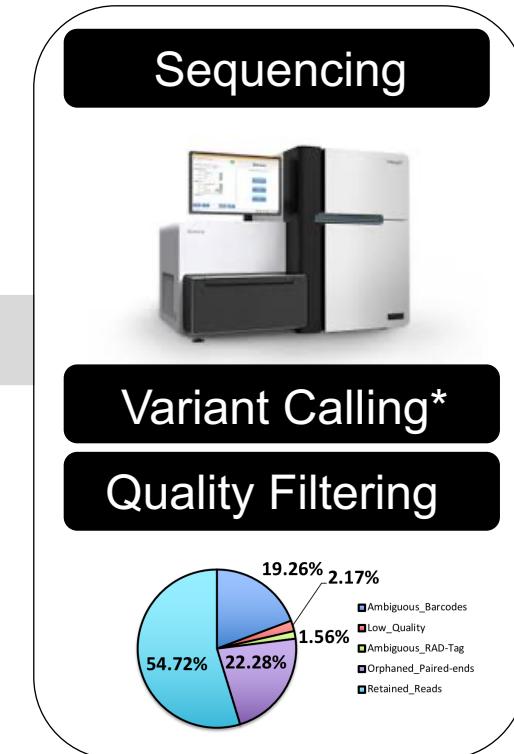
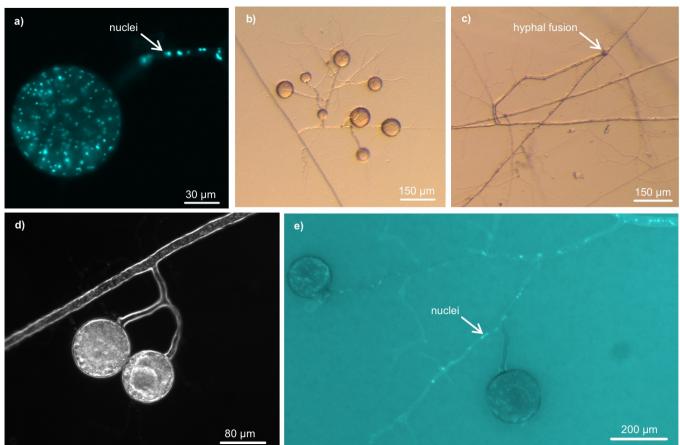


CUSTOM PIPELINES

MY PROJECT

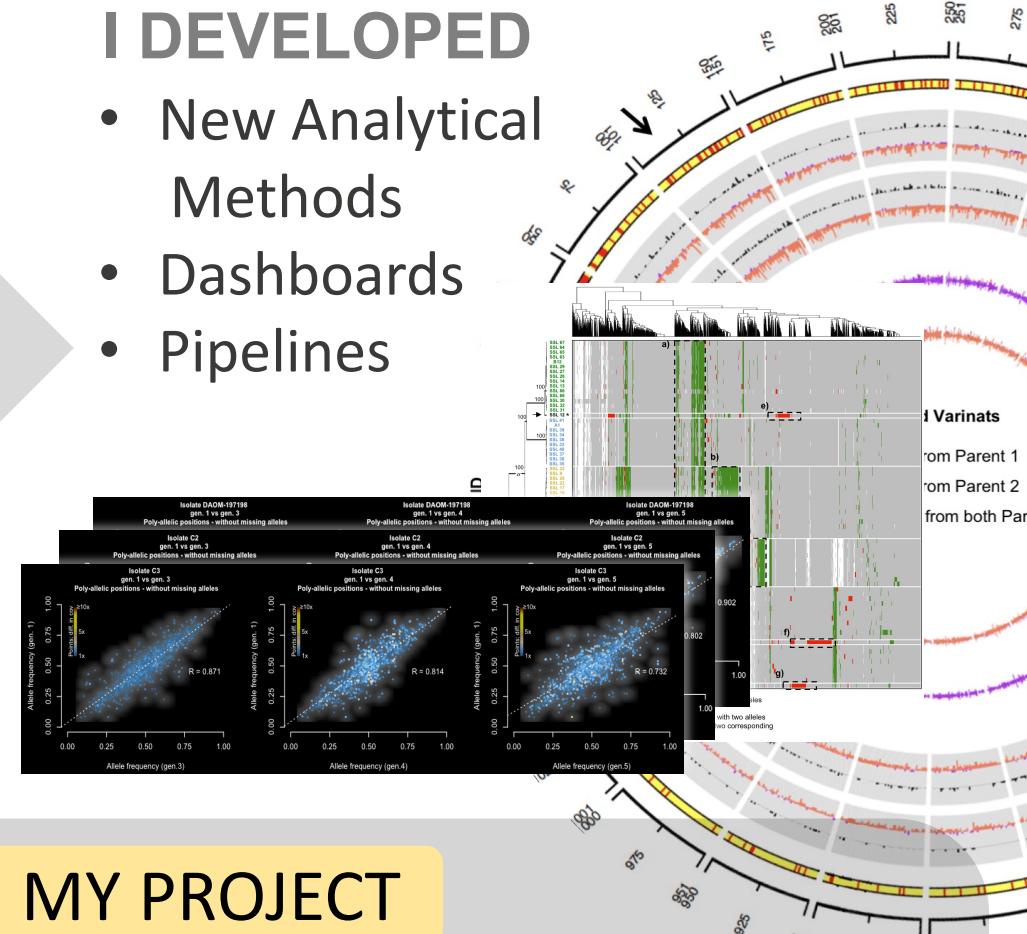
### MY MAIN EXPERIMENT

- Took 5 years
- >25.000 bio-samples
- 3TB of raw data



### I DEVELOPED

- New Analytical Methods
- Dashboards
- Pipelines



### PRODUCTS GENERATED WITH MY PROJECT

IMPACT

**IN VITRO  
CULTURE SYSTEM**

IMPLEMENTED IN  
INDUSTRY

<https://doi.org/10.1007/s00572-017-0763-2>

NGS PROTOCOL

USED OVER 3000x

**60 STRAINS  
OF MICROORGANISMS**

TESTED FOR USE IN  
AGRONOMY

<https://doi.org/10.1371/journal.pone.0226497>

<https://www.biorxiv.org/content/10.1101/830547v1>

I used simulated data and to explain patterns observed in empirical data and to build new pipelines

## EXAMPLE

**One of my experiments could have six possible results, and many possible subtypes**

## MY TASK:

- find relevant features
  - create simulated data corresponding to each scenario
  - Use statistical tests to find best targets
  - Design new cost-efficient experiment to verify it

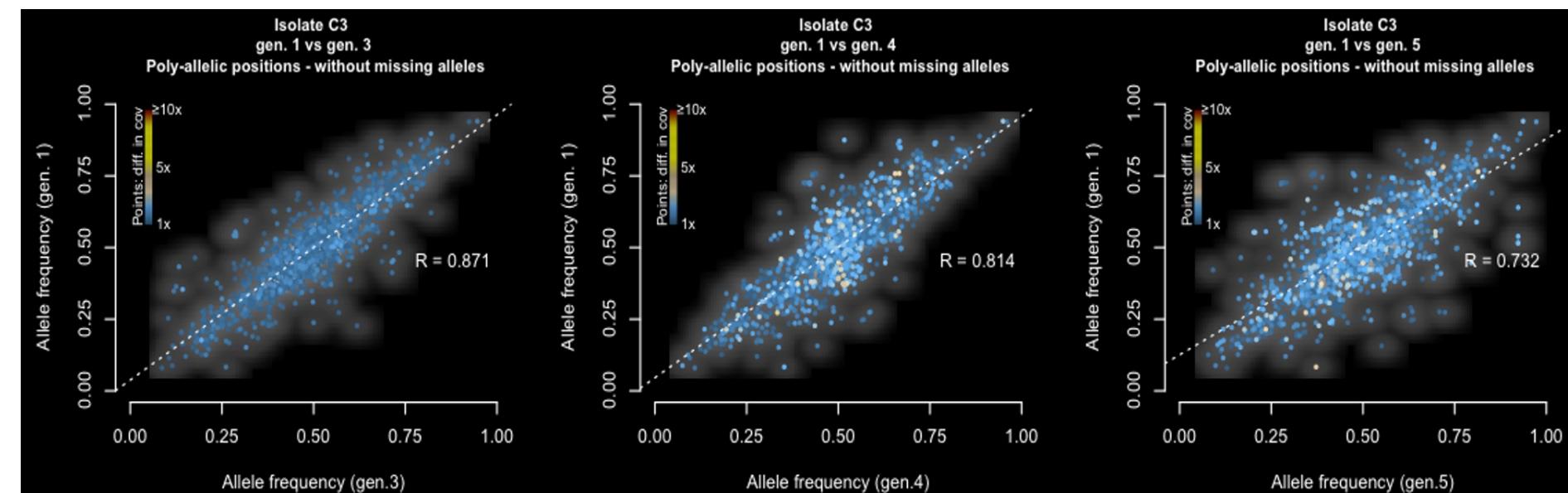
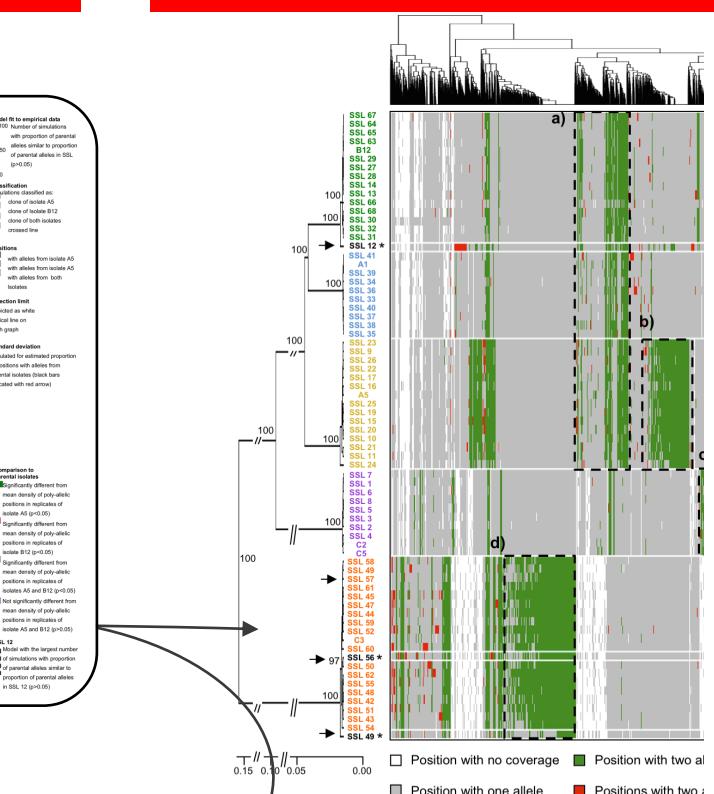
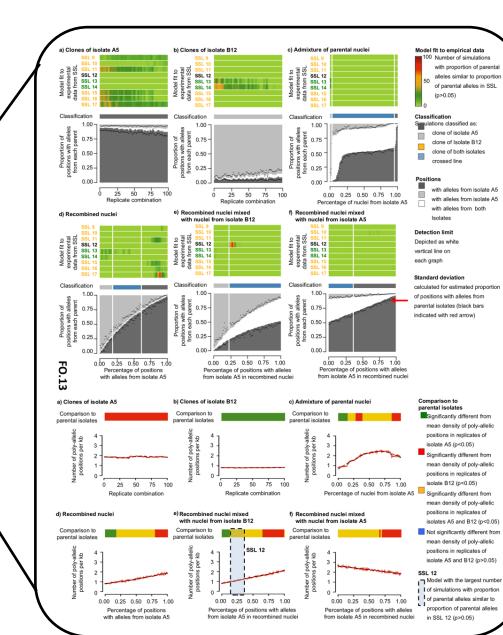
# HOW?

- I used simulated data with known results to build, and test new pipelines
  - I used empirical data as starting point for my simulations
  - I repeated the process over 6 million times for 120.000 features, in thousands of conditions and compared them with empirical data

I created Dashboards to summarize the results

... to see more details if needed

... and how samples relate to each other



In one year  
I created a  
community of over  
2000 AI specialists,  
thanks to  
systematic work and  
connection with  
strong brands

**SwissAI**  
Machine Learning  
Meetup



## Pawel Rosikiewicz

**Founder and Team Leader**  
of SwissAI Machine Learning Meetups  
one of the largest group in Europe

[www.SwissAI.org](http://www.SwissAI.org)

